

This application is submitted in the name of the following inventor(s):

| <i>Inventor</i> | <i>Citizenship</i> | <i>Residence City and State</i> |
|------------------|--------------------|---------------------------------|
| Andrew CONWAY | Australia | Stanford, California |
| Peter EASTMAN | United States | Belmont, California |
| Howard SNORTLAND | United States | Menlo Park, California |
| Barrett EYNON | United States | Menlo Park, California |

The assignee is *Silicon Genetics*, a California corporation having an office at 2601 Spring Street, Redwood City, CA 94063.

TITLE OF THE INVENTION

Autonomous Data Mining

BACKGROUND OF THE INVENTION

I. Field of the Invention

This invention relates to autonomous data mining, such as for example autonomous data mining in large sets of gene expression data.

2. *Related Art*

In computer systems having relatively large amounts of data, such as recorded in a database system or other system for storage and retrieval of data, it is sometimes desirable to review that data to find if there are relationships between data elements that were previously unconfirmed or even unknown. This process is sometimes called “data mining,” and is typically applied to programmed processes that are applied to relatively large databases. For example, searching a large database of stock data for those securities that meet predetermined criteria for capitalization and earnings would be a form of data mining.

Known methods of data mining include “clustering,” that is, attempting to divide the multiple data elements into a relatively small set of clusters. Other known methods include applying statistical methods to best-fit a predetermined relationship against the set of data, so as to determine a set of parameters for the predetermined relationship. These other known methods include multiple linear regression and other statistical and stochastic techniques. While these methods of the known art can generally achieve the purpose of evaluating predetermined relationships against a relatively large set of data, they are of course subject to the drawbacks of all statistical methods, which is that they can only deliver a probabilistic assessment of the predetermined relationship against the set of data.

One problem with the known art is that the researcher or other person (that is, a “user”) must have a predetermined relationship in mind before attempting to apply it against the set of data. For example, when searching a large database of stock data, the user must have a predetermined relationship and a set of predetermined stock parameters in mind for evaluation before known data mining techniques can evaluate whether that predetermined relationship applies well to that set of predetermined stock parameters. This could be referred to as a hypothesis-generating problem.

A second problem with the known art is that the predetermined relationship might have little or no relationship to domain-specific knowledge about the set of data. For example, when searching a large database of stock data, the user might request evaluation of a predetermined relationship among a set of predetermined stock parameters that have, in any real-world model of the stock market, no relationship to each other (such as, for example, whether stocks with a price/earnings ratio that is a prime number occur more frequently when the Moon is in the Aries constellation). This could be referred to as the uninteresting-hypothesis problem.

A third problem with the known art is that the predetermined relationship and the set of data must be determined ahead of the operation of the data mining method. For example, when searching a large database of stock data, the user must assure that all needed data is available before attempting to perform data mining. This could be referred to as the known-database problem.

1 All three of these problems are particularly acute in the field of scientific
2 research into gene expression.

3
4 First, databases of gene expression data have been collected by researchers
5 and are often made available to each other, either in the context of academic research or
6 in the context of pharmaceutical or other for-profit research and development. These da-
7 tabases are relatively large, and are getting substantially larger as time goes by, both due
8 to work by researchers in obtaining new gene expression data and due to improved meth-
9 ods for obtaining that data in greater quantity and at greater speed. As an emergent con-
10 sequence of the rapid growth of databases of gene expression data, it has become ex-
11 tremely difficult for individual researchers to maintain familiarity even with the scope of
12 data available for review.

13
14 Second, gene expression data includes raw data describing measurements
15 of activity for individual strands of mRNA (messenger RNA). These measurements can
16 differ in response to differing times they were taken, differing patients they were taken
17 from, differing clinical samples from one or more patients, differing medical conditions
18 of the one or more patients, differing prescription or other drugs the patients were under
19 the influence of, differing chemical milieus in which the measurements were taken, and
20 many other possible differing conditions. Collection, recording and publication of gene
21 expression data are known in the art of biochemical research. As might be inferred from
22 this description, sets of gene expression data can be extremely complex, having no im-

1 mediate relationships available to the reviewer of the data. Moreover, new sets of gene
2 expression data are generated from time to time, thus increasing the available pool of
3 gene expression data relatively continuously.

4
5 Third, the research community does not always make these sets of gene ex-
6 pression data available immediately upon production. Sometimes individual sets of data
7 are checked for consistency or quality control. Sometimes one or more researchers have
8 a particular predetermined relationship they would like to evaluate (and publish) before
9 allowing other research groups to access those sets of data. As the number and size of
10 sets of gene expression data becomes larger, and as the number of researchers interested
11 in those sets of data becomes larger, the chance that a valuable set of data is not available
12 to one or more researchers interested in that valuable set of data becomes greater.

13
14 Fourth, the particular biological processes that sets of gene expression data
15 reflect are relatively complex. There are relatively large numbers of genes, activation of
16 each of which possibly affects large subsets of other genes, in ways that are presently not
17 well known. (That is why study of gene expression data is called “research.”) Many of
18 these processes are highly nonlinear, that is, a small change in amount of gene expression
19 for a first gene can result in very large changes in amounts of gene expression for one or
20 more downstream sets of genes. Many of these processes have feedback, feed-forward,
21 or other complex topological loops, so that gene expression for a first gene can have
22 multiple different effects on gene expression for both a second gene and for the first gene

1 itself. Even relatively simple examples known as cell cycles can involve relatively long
2 feedback loops, each element of which itself includes a relatively complex set of interac-
3 tions.

4
5 Known methods of examining gene expression data include examining the
6 data “by hand,” that is, by an interested researcher who formulates hypotheses, performs
7 operations on the data to evaluate those hypotheses, and determines if there is sufficient
8 support for those hypotheses to warrant further experiment or even publication of results
9 of the evaluation. While these known methods generally achieve the goals of finding and
10 publishing interesting and useful statements about gene expression to the research world,
11 they are subject to several drawbacks. As noted above, there is a relatively large amount
12 of gene expression data. The amount of such data is rapidly increasing and is not easily
13 subject to efficient or effective search by human researchers. Researchers do not have
14 adequate time to review all the relevant data. Researchers also do not have adequate time
15 to determine all the relevant data in their field, or in related fields. Researchers also often
16 work in close-knit groups and are therefore not always aware of similar work being per-
17 formed by other researchers. Moreover, as noted above, problems in gene expression
18 analysis are relatively complex, and are therefore not easily subject to “by hand” analysis
19 of extensive data.

20
21 Accordingly, it would be desirable to provide a technique in which data
22 mining is performed with regard to a set of data, possibly interesting hypotheses are for-

1 mulated in response thereto, and those hypotheses are reported. In one aspect of the in-
2 vention, this technique can be achieved by performing a robotic process with regard to a
3 set of data in a database, so as to formulate potentially interesting hypotheses and so as to
4 communicate those hypotheses to researchers and other persons having an interest
5 therein.

7 SUMMARY OF THE INVENTION

8
9 The invention provides a method and system for performing data mining
10 autonomously with regard to a set of data, and formulating hypotheses in response
11 thereto. An autonomous software element collects sets of data (such as gene expression
12 data), along with collateral data (such as information about published papers, individual
13 researchers, and known relationships between genes), into a unified but extensible data-
14 base. The autonomous software element formulates possibly interesting hypotheses with
15 regard to data findable in the database. The autonomous software element evaluates each
16 such possibly interesting hypothesis against the data in the database, providing a result
17 that relates the possibly interesting hypothesis against a probability it could have occurred
18 by chance. The autonomous software element rates each possibly interesting hypothesis,
19 in response to multiple factors, such as for example (1) if they relate to genes one or more
20 researchers have indicated they are interested in, (2) if they relate to genes for which
21 there are published papers, (3) if they are relatively simple and relatively unlikely to be
22 due to chance, (4) if they are relatively consistent or relatively inconsistent with domain-

1 specific knowledge about known effects of gene expression. The autonomous software
2 element reports those hypotheses to researchers and other interested persons (collectively,
3 “users”), selecting those users who are most likely to be interested in each specific hy-
4 pothesis and who are most interested in being informed of such discoveries.

5
6 As an emergent consequence of the invention, autonomous data mining of a
7 database including a set of data (such as for example a set of gene expression data) pro-
8 duces a set of possibly interesting hypotheses, each of which can itself be graded with re-
9 gard to a set of selected parameters. Additional data, such as research interests, published
10 papers, and the like, can provide information for grading the hypotheses and for grading
11 communications to researchers. The set of hypotheses and additional data itself provides
12 a database for use in determining which hypotheses to report to which researchers and
13 which other persons. This has the advantages of allowing (1) reporting of interesting hy-
14 potheses to those researchers most likely to take action thereon, and (2) filtering of which
15 hypotheses to report in response to researcher preferences.

16
17 The invention provides an enabling technology for a wide variety of appli-
18 cations for data mining and hypothesis testing, to obtain substantial advantages and capa-
19 bilities that are novel and non-obvious in view of the known art.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 shows a block diagram of a system used with an autonomous data mining tool.

Figures 2A and 2B show a flow diagram of a method of operating a system as in figure 1.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

In the following description, a preferred embodiment of the invention is described with regard to preferred process steps and data structures. Embodiments of the invention can be implemented using general-purpose processors or special purpose processors operating under program control, or other circuits, adapted to particular process steps and data structures described herein. Implementation of the process steps and data structures described herein would not require undue experimentation or further invention.

Related Applications

Inventions described herein can be used in conjunction with inventions described in the following document(s).

- 1 • U.S. Provisional Patent Application Serial No. 60/243,650, filed October 26, 2000,
2 in the name of the same inventor, then attorney docket number 7550-0001 (now
3 attorney docket number 208.1001.01), titled "Enterprise-wide Data Mining."
4

5 Each of these documents is hereby incorporated by reference as if fully set
6 forth herein. This application claims priority of each of these documents. These docu-
7 ments are collectively referred to as the "Incorporated Disclosures."
8

9 *System Elements*

10
11 Figure 1 shows a block diagram of a system used with an autonomous data
12 mining tool.
13

14 A system 100 includes an autonomous software element 110, a local data-
15 base 120, a communication link 130, a set of external databases 140, and a set of users
16 150.
17

18 The autonomous software element 110 includes a database access module
19 111, a hypothesis formulation module 112, a hypothesis evaluation module 113, an inter-
20 est-matching module 114, and a publication module 115.
21

1 In a preferred embodiment, the autonomous software element 110 is dis-
2 posed for execution as an application program under control of operating system software
3 on any general-purpose computer workstations 116, such as a PC, Macintosh, or another
4 type of workstation. The workstation has a processor, program and data memory, mass
5 storage, an input device, and a display device. The workstation 116 can run without su-
6 pervision, or can be controlled by an operator 117. Computer workstations are known in
7 the art of computing. The workstation 116 may include a general-purpose computer
8 workstation, a laptop computer, a handheld or "palmtop" computer, or another type of
9 communication or computing device. Due to the computational and memory require-
10 ments of the preferred embodiment, the workstation 116 preferably includes substantial
11 processor, memory and mass storage resources. In other embodiments, the workstation
12 116 may include a plurality of workstations, so as to share memory and other computa-
13 tional resources.

14
15 In a preferred embodiment, each operator 117 includes one or more human
16 operators of an individual workstation 116. However, in alternative embodiments, one or
17 more operators 117 may include a proxy, such as an artificial intelligence program, a da-
18 tabase-querying program, a web-browsing program, some other form of interface pro-
19 gram to one or more actual human beings or other quasi-intelligent devices, or other con-
20 trol programs capable of generating requests for information or responding to responses
21 to those requests.

1 The database access module 111 includes executable software disposed for
2 execution by the workstation 116, and responsive to commands from the operator 117,
3 capable of communication with the local database 120. The database access module 111
4 communicates with the local database 120 by sending database requests, and receiving
5 database responses, from the local database 120. The database access module 111 ob-
6 tains sets of gene expression data from the local database 120, records them in data mem-
7 ory for the workstation 116, and provides those sets of gene expression data to the hy-
8 pothesis formulation module 112, the hypothesis evaluation module 113, the interest-
9 matching module 114, and the publication module 115.

11 The hypothesis formulation module 112 includes executable software dis-
12 posed for execution by the workstation 116, and responsive to commands from the op-
13 erator 117, capable of receiving gene expression data and operating on that gene expres-
14 sion data (as further described below) so as to generate a set of hypotheses.

16 The hypothesis evaluation module 113 includes executable software dis-
17 posed for execution by the workstation 116, and responsive to commands from the op-
18 erator 117, capable of receiving a set of hypotheses from the hypothesis formulation
19 module 112, and capable of evaluating those hypotheses (as further described below) for
20 relative interest and likelihood.

1 The interest-matching module 114 includes executable software disposed
2 for execution by the workstation 116, and responsive to commands from the operator
3 117, capable of receiving a set of hypotheses from the hypothesis evaluation module 113,
4 and capable of matching those hypotheses (as further described below) with a set of data
5 regarding users, so as to evaluate which users are most likely to be interested in those hy-
6 potheses, and so as to evaluate which users are most likely to want to be informed
7 thereof.

8
9 The publication module 115 includes executable software disposed for exe-
10 cution by the workstation 116, and responsive to commands from the operator 117, capa-
11 ble of receiving a set of hypotheses and a set of corresponding users from the interest-
12 matching module 114, and capable of publishing those hypotheses to those users. In a
13 preferred embodiment, the publication module 115 generates an email message for each
14 such hypothesis to each such user, and sends that email message using the communica-
15 tion link 130 to the target user.

16
17 The local database 120 includes a database system 121, having a set of gene
18 expression data elements 122 and a set of collateral data elements 123 recorded therein.

19
20 The database system 121 is accessible to the autonomous software element
21 110 by relatively local techniques. In a preferred embodiment, the local database 120 is
22 coupled directly to the workstation 116. In alternative embodiments, the local database

120 may be accessible to the workstation 116 using a LAN (local area network), using the communication link 130, or using any other technique by which the autonomous software element 110 can relatively rapidly and reliably access relatively large amounts of data recorded in the local database 120. In a preferred embodiment, the local database 120 includes a standard relational or object-oriented database system, such as the "Oracle" product available from Oracle Corporation of Redwood City, California.

The gene expression data elements 122 each include information regarding a particular gene expression measurement, including information regarding times they were taken, patients they were taken from, clinical samples from one or more patients, medical conditions of the one or more patients, prescription or other drugs the patients were under the influence of, the chemical milieu in which the measurements were taken, and other relevant conditions.

The collateral data elements 123 each include information such as the following:

- information regarding particular genes, possibly including information regarding known or hypothesized interactions with other genes, gene sequence data and organism sequence data.

- 1 • particular researchers, possibly including information regarding (1) individual
2 genes or gene expressions those researchers have particular interest in, (2) the de-
3 gree of receptivity to messages from the autonomous software element those re-
4 searchers have;
- 5
- 6 • published results and/or papers, possibly including (1) further domain-specific
7 knowledge regarding individual genes or gene expressions, (2) known relation-
8 ships between individual researchers, such as having collaborated on papers to-
9 gether or being located at the same or related institutions.
- 10

11 The communication link 130 includes a physical communication link, con-
12 trol software, and communication protocols for devices attached thereto to send and re-
13 ceive messages. The communication link 130 can include an internet, intranet, extranet,
14 VPN (virtual private network), private or public switched network, ATM network, LAN
15 (local area network), WAN (wide area network), direct communication line, shared
16 memory communication, or any other technique capable of performing the functions de-
17 scribed herein. In a preferred embodiment, the communication link 130 includes an
18 internet and a LAN connection to the internet.

19

20 The external databases 140, similar to the local database 120, include data-
21 base systems 141, each having a set of data elements 142 recorded therein.

22

1 The database systems 141 are accessible to the autonomous software ele-
2 ment 110 using the communication link 130. The external databases 140 may include
3 standard relational or object-oriented database systems, or may include text files or other
4 data formats readable by the autonomous software element 110 (including HTML).

5
6 The gene expression data elements 142 in the external databases 140 in-
7 clude sets of gene expression data. In a preferred embodiment, gene expression data
8 elements 142 include elements similar to the gene expression data elements 122.

9
10 The collateral data elements 143 in the external databases 140 include in-
11 formation regarding published papers, individual researchers, and known relationships
12 between genes. In a preferred embodiment, collateral data elements 143 include elements
13 similar to the collateral data elements 123.

14
15 The users 150 each include researchers, other interested persons, or inter-
16 ested groups or institutions, each having a workstation 116 and an operator 117, similar
17 to the autonomous software element 110.

18
19 *Method of Operation*

20
21 Figures 2A and 2B show a flow diagram of a method of operating a system
22 as in figure 1.

1 A method 200 includes a set of flow points and process steps as described
2 herein.

3
4 Although by the nature of textual description, the flow points and process
5 steps are described sequentially, there is no particular requirement that the flow points or
6 process steps must be sequential. Rather, in various embodiments of the invention, the
7 described flow points and process steps can be performed in a parallel or pipelined man-
8 ner, either by one device performing multitasking or multithreading, or by a plurality of
9 devices operating in a cooperative manner. Parallel and pipelined operations are known
10 in the art of computer science.

11
12 At a flow point 210, the autonomous software element 110 is ready to be-
13 gin.

14
15 At a step 211, the database access module 111 retrieves a set of gene ex-
16 pression data elements 122 from the external databases 140 and records them in a unified
17 but extensible local database 120. In a preferred embodiment, the database access mod-
18 ule 111 operates from time to time to re-perform this step, even if other parts of the
19 autonomous software element 110 are operating in parallel.

At a step 212, the hypothesis formulation module 112 formulates a possibly interesting hypothesis. To perform this step, the hypothesis formulation module 112 does one or more of the following:

- Comparing sets of genes: The hypothesis formulation module 112 identifies a test set G_T of genes, as follows: (1) The hypothesis formulation module 112 prepares a set of genes G_C found by clustering, as further described below. (2) The hypothesis formulation module 112 compares the set of genes G_C found by clustering with other known sets of genes G_K found by other means. For example, known sets of genes G_K found by other means can include genes found by matching keywords in the external databases 140, or other private or public databases. (3) The hypothesis formulation module 112 performs a statistical technique to determine a probability of overlap between the clustering set G_C and the known set G_K . If the statistical technique indicates that the probability is relatively low, the hypothesis formulation module 112 concludes that the clustering list G_C would be a useful test set G_T .
- Comparing upstream sequences: (1) The hypothesis formulation module 112 prepares a set of genes G_C found by clustering, using a clustering technique as further described for “Comparing sets of genes.” (2) The hypothesis formulation module 112 gathers from database 120 sequences of nucleotides upstream from the genes in the clustering set G_C . (3) The hypothesis formulation module 112 determines if there are relatively short, such as 5-20 bases, sequences that are more common up-

stream of the clustering set G_C than of other genes in the genome. This could be an interesting observation as it could indicate that these sequences may be regulatory elements for a set of co-regulated genes.

- Quality control / data collection: (1) The hypothesis formulation module 112 examines those gene expression data elements 122 from pairs of similar experiments, to determine if their data was poorly collected, taking into account the technology used. For example, for 2 color experiments, hypothesis formulation module 112 checks whether genes with high expression in one experiment tend to be significantly over-expressed or under expressed in the other experiment.
- Quality control / experiment replication: (1) The hypothesis formulation module 112 examines those gene expression data elements 122 from replicated experiments, that is, the same experiment performed at differing times. (2) The hypothesis formulation module 112 determines if there is any statistically significant difference between the replicated experiments. For example, the hypothesis formulation module 112 can compute if the number of genes for each experiment that contain extreme values indicates that the numbers are unevenly distributed. A low probability indicates a possibility that data has been incorrectly measured or entered

- 1 • Pathway completion: (1) The hypothesis formulation module 112 examines
2 “pathways” of genes, such that one or more genes, in turn activates the next. A
3 pathway can be represented visually as a network of genes with relationships de-
4 fined between the genes either of connectedness or physical layout in a 1D, 2D or
5 higher dimensional graph. (2) The hypothesis formulation module 112 looks for
6 sets of genes G_C that might fit into gaps in the pathway, such as by examining
7 each possible gene for each gap in response to (a) a distance metric of expression
8 values for the possible gene and the genes known to be on the pathway, or (b) a
9 distance metric of physical distance for the possible gene and the genes known to
10 be on the pathway. If the relative distances between genes determined with re-
11 spect to expression data are similar in pattern to the relative distances between
12 genes determined with respect to the pathway diagram, then that gene may fit in
13 that place on the diagram. This is useful as it can place previously unknown genes
14 on a pathway.
15
16 • Parameter variation: The hypothesis formulation module 112 looks for genes that
17 behave one way in some experiments, and a different way in similar experiments
18 varying only by one or a few parameters. Statistical validation of this can be
19 achieved by estimating experimental error for each point based upon replicates
20 and the absolute expression values.
21

- 1 • Unusual behavior: The hypothesis formulation module 112 looks for genes that
2 behave similarly in most experiments, but significantly differently in a few ex-
3 periments. This is achieved by looking at correlation values in the different ex-
4 periments, and building a model for the correlations of different genes. If this cor-
5 relation is significantly different in a small number of experiments, this will be
6 considered important.
- 7
- 8 • Recent data: The hypothesis formulation module 112 looks for experiments or
9 annotations that have been recently added to the local database 120.

10

11 As part of this step, the “Comparing sets of genes” and “Comparing up-
12 stream sequences” techniques (and possibly also the “Parameter variation” and “Unusual
13 behavior” techniques) use clustering to select a set of genes G_C . There are several known
14 techniques for clustering known in the art of data mining.

15

16 To perform clustering, the hypothesis formulation module 112 selects a set
17 E_C of gene expression data elements 122 on which to perform clustering. In a preferred
18 embodiment, the set of gene expression data elements 122 is responsive to a common set
19 of information. For example, the common set of information can be one or more of: (1)
20 the gene expression data elements 122 all relate to the same patient, the same drug test, or
21 the same chemical environment; (2) the gene expression data elements 122 all relate to
22 the same keyword; (3) the gene expression data elements 122 all relate to the same re-

searcher or the same time period; (4) the gene expression data elements 122 all relate to interaction of the same first gene G_1 with a set of other genes.

In a preferred embodiment, the clustering technique uses hierarchical clustering using nested subtrees of partitions. This is a known technique in the art of statistics.

In alternative embodiments, the clustering technique may include k-means clustering (this technique is known in the art of statistics). In k-means clustering, the set E_C of gene expression data elements 122 is divided into a pre-selected number of clusters E_1, E_2, \dots, E_N . Each cluster has a cluster center, initially chosen at random. Each one of the gene expression data elements 122 is assigned to one of the clusters in response to its distance from the cluster center (the one with minimum "distance," according to a selected distance metric). Cluster centers are moved in response to those records assigned to each cluster. This process is repeated until the cluster centers are static to a selected degree.

In other alternative embodiments, the clustering technique may include any other effective clustering technique, such as the following: self-organized maps or user-directed clustering "by hand."

1 At a step 213, the hypothesis evaluation module 113 tests the hypothesis.

2 To perform this step, the hypothesis evaluation module 113 performs the following sub-
3 steps:

- 4
- 5 • At a sub-step 213(a), the hypothesis evaluation module 113 retrieves gene expres-
6 sion data elements 122 relevant to the hypothesis.
 - 7
 - 8 • At a sub-step 213(b), the hypothesis evaluation module 113 performs statistical
9 tests on the gene expression data elements 122, to determine whether it is possible
10 to confidently reject the possibility that the hypothesis was true by chance.
 - 11
 - 12 • At a sub-step 213(c), if the statistical tests indicate that the hypothesis evaluation
13 module 113 can confidently reject the possibility that the hypothesis was true by
14 chance, the hypothesis is marked as publishable because it is possibly interesting.
 - 15

16 At a step 214, the database access module 111 retrieves a set of collateral
17 data elements 123 from the external databases 140 and records them in a unified but ex-
18 tensible local database 120. In a preferred embodiment, the database access module 111
19 operates from time to time to re-perform this step, even if other parts of the autonomous
20 software element 110 are operating in parallel.

1 At a step 215, the interest-matching module 114 compares the publishable
2 hypothesis with collateral data elements 123, so as to determine to which users 150 to
3 publish the publishable hypothesis. In a preferred embodiment, this step is performed in
4 parallel or asynchronously with the other steps of the method 200, as it is possible for a
5 user 150 to change their interests while the autonomous software element 110 is in op-
6 eration. To perform this step, the interest-matching module 114 performs the following
7 sub-steps:

- 8
- 9 • At a sub-step 215(a), the interest-matching module 114 retrieves collateral data
10 elements 123 for each user 150.
 - 11
 - 12 • At a sub-step 215(b), the interest-matching module 114 determines a interest
13 ranking for the publishable hypothesis for each user 150. In a preferred embodi-
14 ment, the interest ranking is responsive to (1) expressions of interest or disinterest
15 by each user 150, (2) general factors tending to indicate interest or disinterest by
16 all users 150, and (3) reliability of data underlying the publishable hypothesis.

17 These collectively include the following factors, amongst others:

- 18
- 19 ○ whether the user 150 expressed an interest in one of the experiments, genes
20 or gene lists in the publishable hypothesis;

- whether the user 150 was an author of a paper, or a provider of data for, any part of the publishable hypothesis;
- whether the user 150 expressed an interest in any key word associated with one of the experiments, genes or gene lists in the publishable hypothesis;
- whether the user 150 expressed an interest in a “type” of hypothesis;
- whether the user 150 expressed an interest in any author of a paper, or a provider of data for, any part of the publishable hypothesis;
- whether the user 150 expressed an interest in the genome the genes or gene lists are part of, for the genes or gene lists in the publishable hypothesis;
- a measure of age of the gene expression data elements 122 used to form or test the publishable hypothesis;
- a measure of time since the gene expression data elements 122 used to form or test the publishable hypothesis were themselves published;
- a measure of time since the result was generated by the system modules 112 and 113;

- a measure of a number of users 150 who are considered to be interested in the publishable hypothesis or to whom the hypothesis has already been sent;
- a measure of a number of users 150 who have responded with feedback regarding whether they were interested in the publishable hypothesis;
- a measure of similarity between the publishable hypothesis and other hypotheses other users 150 have considered interesting;
- a measure of how many publishable hypotheses, or other knowledge, already exist regarding the genes or gene lists in the publishable hypothesis.
- At a sub-step 215(c), if the statistical tests indicate that the hypothesis evaluation module 113 can confidently reject the possibility that the hypothesis was true by chance, the hypothesis is marked as publishable and possibly “interesting.”

At a step 216, the publication module 115 generates a publication regarding the publishable hypothesis. The publication can include a database or other data file in a specified format, an HTML (hypertext markup language) page, or an email message. In a preferred embodiment, the publication module 115 generates an HTML page including the largest portion of information associated with the publishable hypothesis, and sends

1 an email message with a synopsis (including a pointer to that HTML page) to each inter-
2 ested user 150.

3
4 At a flow point 220, the autonomous software element 110 has completed
5 one cycle of finding and publishing a publishable hypothesis. The method 200 continues
6 with the flow point 210, unless interrupted by the operator 117.

7
8 *Generality of the Invention*

9
10 The invention has general applicability to various fields of use, not neces-
11 sarily related to the particular data, databases, data sets, or uses described above. For ex-
12 ample, these fields of use can include one or more of, or some combination of, the fol-
13 lowing:

- 14
15 • experimental data or processes in scientific fields other than biochemistry, such as
16 flow dynamics, materials science, radiochemistry, weather; and
17
18 • experimental data or processes in non-“scientific” fields, such as financial instru-
19 ments, operations research, product marketing, risk analysis.
20
21

1 Other and further applications of the invention in its most general form,
2 will be clear to those skilled in the art after perusal of this application, and are within the
3 scope and spirit of the invention. Although preferred embodiments are disclosed herein,
4 many variations are possible which remain within the concept, scope, and spirit of the in-
5 vention, and these variations would become clear to those skilled in the art after perusal
6 of this application.

7
8 *Technical Appendix*

9
10 This application includes a technical appendix, hereby incorporated by ref-
11 erence as if fully set forth herein. The technical appendix includes the following:

- 12
13 • an example report generated and published by the automated software element.

14
15 Although the technical appendix relates to a particular preferred embodi-
16 ment of the invention, the information in the technical appendix is also applicable to al-
17 ternative embodiments of the invention, and should be read to indicate the possible scope
18 and spirit of the invention. The technical appendix is not intended to be limiting of the
19 scope or spirit of the invention in any way.